

Operator scaling and matrix sensing

Plan:

- Matrix & operator scaling recap.
- Regression, mirror descent, Matrix scaling
- Sensing, operator scaling, moment maps.
- Open directions

Overparametrized linear regression + matrix sensing.

- Sparse regression:

$u: \mathbb{R}^N \rightarrow \mathbb{R}^M$ linear map, $b \in \mathbb{R}^M$

- Want to solve $u(A) = b$ for $A \in \mathbb{R}_{\geq 0}^N$.

- If $N \gg M$, many solutions;
want a sparse one.

- Proxy for sparsity: $\|A\|_1$.

- Minimize quadratic loss $\|u(A) - b\|^2$

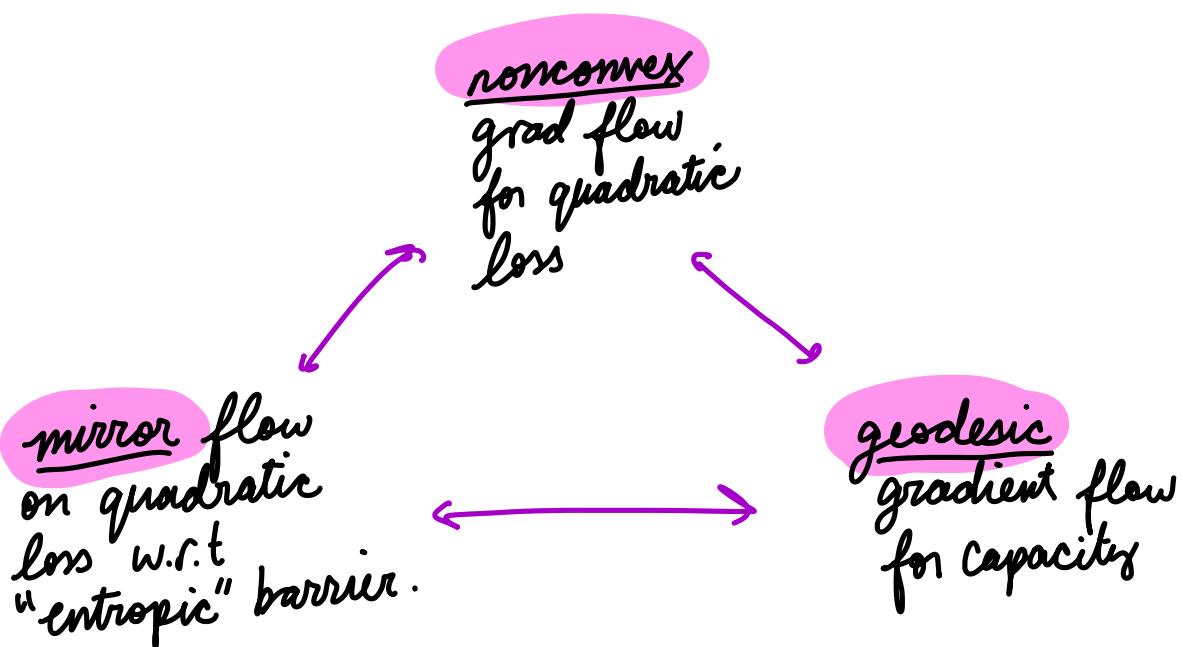
- Use nonconvex reparametrization

$$A = X \circ \underset{\text{coordinatewise}}{\underset{\uparrow}{X}} \text{ product.}$$

- Matrix sensing: Similar thing;
 $u: \text{Mat}(N \times N) \rightarrow \mathbb{R}^M$
want to solve $u(p) = b$
for $p \succcurlyeq 0$.
- In that case, reparametrize $p = XX^T$.
for $X \in \text{Mat}(N, k)$.

Three flows

- For regression case, identical:



- More detail to follow.
- Main Question: also true for matrix sensing?
- Interesting even for special case called operator scaling.
- Example: Matrix scaling

Given $A \geq 0 \in \text{Mat}(m, n)$,
 $(r, c) \in \mathbb{R}^m \times \mathbb{R}^n$,

find diagonal matrices D_1, D_2 s.t.

$\tilde{A} = D_1 A D_2$ has

$$\tilde{A} \mathbb{I} = r, \quad \tilde{A}^\top \mathbb{I} = c.$$

- Rephrase: Let

$$\triangleright u(A) := (A\mathbb{1}, A^T \mathbb{1})$$

$$\triangleright b := (r, c)$$

- So we are looking to solve $u(A) = b$ a very special way.

- Note $u: \mathbb{R}^{mn} \rightarrow \mathbb{R}^m \times \mathbb{R}^n$; $N = mn$, $M = m+n$

Mirror vs. nonconvex flow.

- Back to general $u: \mathbb{R}^N \rightarrow \mathbb{R}^M$.

- write $A = X \circ X$, let $X(t)$

gradient flow of $q(X) := \|u(X \circ X) - b\|_2^2$.

- To minimize $\|A\|_1$,

Start at $X(0) = \alpha \mathbb{1}$ for $\alpha \rightarrow 0$.

Fact: Let $A(\infty) = \lim_{t \rightarrow \infty} X(t)^\dagger X(t)$; then

$$A(\infty) \xrightarrow[\alpha \rightarrow 0]{} \underset{u(A)=b}{\operatorname{argmin}} \|Ax\|_1.$$

(Gunasekar et.al. 2017)

- Gunasekar et. al. conjectured
Same for matrix sensing; recently
disproven! (Li, Luo, Lyu 2020)
- Example is 4×4 ; converges to
lowest rank.

Stronger Fact: For any $A(0) \geq 0$,

$$A(\infty) = \underset{u(A) = b}{\operatorname{argmin}} D_{KL}(A \parallel A(0))$$

** whenever min is finite.*

where

$$D_{KL}(A \parallel B) = \sum B_i - \sum A_i + \sum A_i \log \frac{A_i}{B_i}.$$

Proof Sketch:

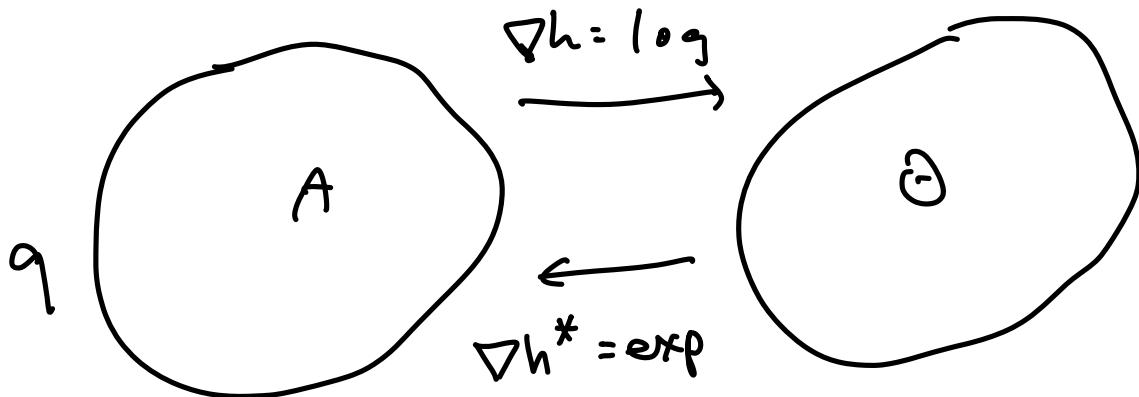
- suffices to show $A(t)$

mirror flow of loss q w.r.t. barrier

$$h(A) = \sum A_i (\log A_i - 1)$$

"negative entropy".

b/c D_{KL} is Bregman divergence w.r.t. h .



- flow: $\dot{\theta} = -\nabla_A f(A) \Big|_{A=\exp \theta}$
 $A(t) := \exp \Theta(t),$

- both ODE's for A are

$$\dot{A}(t) = -\mu^* (\mu(A) - b) \circ A.$$

- Here $\mu^*: \mathbb{R}^M \rightarrow \mathbb{R}^N$ is adjoint of μ .
- Example: for matrix scaling,

$$\mu^*(y, z) = y \mathbf{1}^T + \mathbf{1} z^T. \quad \square$$

Geodesic flow for matrix scaling

- capacity $\text{cap}_{r,c}(A)$ defined by

$$\inf_{\substack{y \in \mathbb{R}^m \\ z \in \mathbb{R}^n}} \sum_{ij} A_{ij} e^{y_i + z_j} - y \cdot r - z \cdot c$$

$\curvearrowright := f(y, z)$

- Sinkhorn is just alt. min;
scaling is $e^{\text{diag}(y)} A e^{\text{diag}(z)}$.

- Consider gradient flow for $f(x, y)$.

$$A(t) := e^{\text{diag } y(t)} A e^{\text{diag } z(t)}$$

- Then $\dot{y}(t) = F - A(t)^T \mathbf{1}$
 $\dot{z}(t) = C - A(t)^T \mathbf{1}$.

- Or with $b = (F, C)$, $\mu(A) = (A\mathbf{1}, A^T\mathbf{1})$,

$$(\dot{y}, \dot{z}) = b - \mu(A).$$

- Now
- $$\begin{aligned}\dot{A}(t) &= \text{diag}(\dot{y}) A + A \text{diag}(\dot{z}) \\ &= (\dot{y} \mathbf{1}^T + \mathbf{1} \dot{z}^T) \circ A \\ &= \mu^*(\dot{y}, \dot{z}) \circ A\end{aligned}$$

$$= -\mu^*(\mu(A) - b) \circ A .$$

- Same as mirror & nonconvex!

- Nothing special about matrix scaling; generalizes for

$$f(\gamma) = \sum_i e^{\mu^*(\gamma)}; A_i - \gamma \circ b,$$

$$A(t) = e^{\mu^*(\gamma(t))} \circ A$$

- check $e^{\mu^*(y,z)} \circ A = e^{\text{diag}(y)} A e^{\text{diag}(z)}$
for matrix scaling.
 \Rightarrow flow stays in group orbit!

Operator Scaling

Matrix Scaling

Nonnegative $n \times m$ matrix A

"(unnormalized) joint distribution" on $[m] \times [n]$

Alice A Bob

Correspondence: $A \mapsto \rho := \sum_{ij} A_{ij} (e_i \otimes e_j) (e_i^* \otimes e_j^*)$

Marginals:

$$\mu(A) = (A^T \mathbb{1}, A \mathbb{1})$$

$$E \in \mathbb{R}_{\geq 0}^m \times \mathbb{R}_{\geq 0}^n$$

Operator Scaling

$m \times m$ PSD matrix ρ .

"Unnormalized" joint quantum state on $\mathbb{C}^m \otimes \mathbb{C}^n$.

Alice ρ Bob

$$\mu(\rho) = (\text{Tr}_2 \rho, \text{Tr}_1 \rho)$$

$$\mu(\rho) \in \text{PSD}(m) \times \text{PSD}(n)$$

$$(\text{Tr}_2 \rho)_{ij} = \text{Tr}(E_{ij} \otimes I_n \rho)$$

$$(\text{Tr}_1 \rho)_{ij} = \text{Tr}(I_m \otimes E_{ij} \rho)$$

Scaling: $D_1 \succcurlyeq 0, D_2 \succcurlyeq 0$
diagonal

$$(D_1, D_2) \cdot A := D_1 A D_2.$$
$$(g_1, g_2) P$$
$$:= g_1 \otimes g_2 \quad (g_1^+ \otimes g_2^+).$$

Problem (operator scaling)

Given $P \succcurlyeq 0$ $m \times m$,

target mary. $R \succcurlyeq 0$ $m \times m$ $C \succcurlyeq 0$ $n \times n$,

find scaling \tilde{P} s.t.

$$\text{Tr}_2 \tilde{P} = R, \quad \text{Tr}_1 \tilde{P} = C.$$

- Why? Null-Gau, MLE's,
Quantum state estimation (quantum OT ??).

Nonconvex v.s. geodesic for operator scaling.

- Nonconvex: write $P = XX^\top$

$$\text{Now } q(X) = \|u(X) - b\|_F^2.$$

- In general for matrix sensing,

$$\dot{X} = -u^*(u(X) - b) X. \quad \text{ignoring constants.}$$

- For operator scaling, $N = n$, $M = m^2 + n^2$.

$$u : \text{Mat}(mn, mn) \rightarrow \text{Mat}(m) \times \text{Mat}(n)$$

$$u^*(Y, Z) = Y \otimes I_n + I_m \otimes Z.$$

- μ^* is a Lie algebra action of $gl_n \times gl_m$ on $\text{Mat}(nm, K)$.
- In particular, μ^* is differential of Lie group action

$$(g, h) \cdot X = (g \otimes h)X .$$

of $GL_n \times GL_m$ on $\text{Mat}(nm, K)$.

i.e. $e^{\mu^*(x, z)} = e^x \otimes e^z .$

- Corollary: For μ of operator scaling,

$$x(t) = g(t) \times h(t) X \quad \forall t \geq 0 ,$$

Hence $p(t) = x(t) x(t)^T$ is a

Scaling of $P = P(0)$.

Geodesic flow:

- Capacity simpler for $R = C = I_n$ case, i.e. $b = (I_n, I_n)$.

$$\text{Cap}(P) := \inf_{\substack{Y \succ 0 \\ Z \succ 0}} \text{Tr} Y \otimes Z P - \log \det Y - \log \det Z.$$

$f(Y, Z).$

- Write $Y = g^t g$, $Z = h^t h$; geodesics thru Y are $Y(t) = g^t e^{Ht} g$ H hermitian.
- Riemannian gradient of $f(Y, Z)$ is then

$$\langle \nabla f, (H_1, H_2) \rangle_f = \partial_{t=0} f(g^t e^{H_1 t} g, h^t e^{H_2 t} h)$$

end up with

$$\nabla_1 f = u_1(\tilde{p}) - I$$

$$\nabla_2 f = u_2(\tilde{p}) - I$$

for $\tilde{p} = g \otimes h P g^+ h^+$.

- Define

$$\dot{g}(t) = -\frac{\nabla f_1(x, z)}{2} g(t)$$

$$\dot{h}(t) = -\frac{\nabla f_2(x, z)}{2} h(t),$$

(need this def because g not det. by X).

$$\Rightarrow \dot{X}(t) = g(t) \otimes h(t) X,$$

hence $\dot{X} = -\frac{1}{2} \mu^* (\mu(X X^T) - b) X.$

- Moral: same flow from geodesic capacity min & nonconvex matrix sensing flow.
- First need μ^* to be Lie algebra representation!
- Doesn't converge to lowest rank solution - if $\rho(0) = A \otimes B$,
 $\rho(t) = A(t) \otimes B(t)$. At.
 $\mu(\rho) = (A^T B, B^T A)$
 $\Rightarrow \rho(t) \rightarrow I \otimes I$; full rank n^2
but always \exists rank 1 solution.
- For 3-tensor scaling w/ $\text{rk } X = 1$
it can take very long to converge

Question: Can we also interpret
as a mirror flow?

should tell us where it converges
in terms of some Bregman div.

- What about matrix sensing NOT coming from lie groups?

$$f(z) = \text{Tr } e^{u^*(z)} p - b \cdot z ?$$

Problem: e^z makes no sense if
 z not matrix or Lie gp. elt.

- What about $R \neq I$, $C \neq I$?

Other R,C

- gets more complicated.

- Want something like

$$\text{cap}_{R,C}(P) = \inf_{Y, Z} \text{Tr} Y \otimes Z P - \text{Tr} R \log Y - \text{Tr} C \log Z,$$

but the last terms are bad.

- We'd like the geodesic gradients to satisfy

$$\text{eq. } \nabla_Y \text{Tr} R \log Y = R;$$

unfortunately not the case.

- But there is a nice q -concave function

$$Y \mapsto \log \det(Y, R)$$

which satisfies

- The catch: requires specific choice of coords. for tangent space, namely

$$Y = g^+ g, Z = h^+ h, \text{ for } g, h \text{ lower tri.}$$

- ② Always exists; Cholesky decomp!

- Set

$$f(Y, Z) = \text{Tr } Y \otimes Z P - \log \det(Y, R) - \log \det(Z, C)$$

- Now geodesic grad (taken this way)
gives

$$\mu(\tilde{P}) - (R, C),$$

which is exactly what we need...

- Why not e.g.

$$\dot{g} = \frac{1}{2} (\mu_1(\tilde{P}) - R) g$$

like before? Not lower tri!

- Must project it to lower tri

$$\hat{g} = \Pi \frac{1}{2} (u, (\tilde{p}) - R) g$$

- Thus $X(t)$ from nonconvex & $g(t) \otimes h(t) X$ from geodesic slightly different.
- Experimentally they converge to same p .
- Any fix? Interestingly even for full observation ($m=1, n>1$).